

WHAT IS CLAIMED IS:

- 1 1. A computing unit in a computing machine, wherein the computing
2 machine performs a plurality of computing operations using the computing unit, the
3 computing unit comprising:
4 a hardware structure that implements networked nodes that receive an input
5 signal and map the input signal to an output signal, wherein nodes in the networked nodes are
6 related by a network of connections between the nodes;
7 a weight matrix input that receives a weight matrix, wherein the weight matrix
8 comprises weights corresponding to the connections; and
9 an activation function input that receives an activation function, wherein the
10 activation function specifies a function for the nodes in the network of nodes,
11 wherein the weight matrix and activation function correspond to a computing
12 operation, wherein the hardware structure maps the input signal through the network of
13 connections in the networked nodes using the corresponding weights of the weight matrix for
14 the connections and the function of the activation function to generate the output signal, the
15 output signal being a result of the computing operation that is determined by the weight
16 matrix and activation function.
- 1 2. The computing unit of claim 1, wherein the networked nodes are
2 arranged in a plurality of layers.
- 1 3. The computing unit of claim 1, wherein the networked nodes form a
2 multi-layer perceptron network.
- 1 4. The computing unit of claim 1, wherein the weight matrix comprises a
2 plurality of sub-matrices.
- 1 5. The computing unit of claim 1, wherein the function of the activation
2 function comprises a linear function.
- 1 6. The computing unit of claim 5, wherein the linear function comprises a
2 unity gain function.
- 1 7. The computing unit of claim 1, wherein the function of the activation
2 function comprises a nonlinear function.

1 8. The computing unit of claim 7, wherein the nonlinear function
2 comprises a sigmoid function.

1 9. The computing unit of claim 7, wherein the nonlinear function
2 comprises a limiter function.

1 10. The computing unit of claim 1, wherein the computing machine
2 comprises an integrated circuit.

1 11. The computing unit of claim 1, wherein the hardware structure
2 comprises one or more units capable of performing multiplication and accumulation
3 operations and one or more activation function units.

1 12. A computing unit in a computing machine, wherein the computing
2 machine performs a plurality of computing operations using the computing unit, the
3 computing unit comprising:

4 an input layer of nodes for receiving an input signal;

5 a middle layer of nodes coupled to the input layer of nodes, wherein the
6 middle layer of nodes are related to the input layer of nodes through a first network of
7 connections, the middle layer configured to process the input signal using middle layer
8 weights corresponding to the first network of connections and an activation function to
9 generate a middle layer signal; and

10 an output layer of nodes coupled to the middle layer of nodes, wherein the
11 output layer of nodes are related to the middle layer of nodes through a second network of
12 connections, the output layer configured to process the middle layer signal using output layer
13 weights corresponding to the second network of connections and the activation function to
14 generate an output signal, the output signal being a result of a computing operation
15 corresponding to the middle and output layer weights and the activation function.

1 13. The computing unit of claim 12, wherein the input, middle, and output
2 layers are constructed into a multi-layer perceptron network.

1 14. The computing unit of claim 12, wherein the input layer of nodes is a
2 multiplexer.

1 15. The computing unit of claim 12, wherein a node in the middle layer of
2 nodes comprises one or more units capable of performing multiply and accumulate operations
3 and one or more activation function units.

1 16. The computing unit of claim 15, wherein one or more units capable of
2 performing multiply and accumulate operations comprise multiply-accumulate units.

1 17. The computing unit of claim 12, wherein a node in the output layer of
2 nodes comprises one or more units capable of performing multiplication and accumulation
3 operations and one or more activation function units.

1 18. The computing unit of claim 17, wherein one or more units capable of
2 performing multiplication and accumulation operations comprise multiply-accumulate units.

1 19. The computing unit of claim 12, wherein the activation function
2 comprises a linear function.

1 20. The computing unit of claim 19, wherein the linear function comprises
2 a unity gain function.

1 21. The computing unit of claim 12, wherein the activation function
2 comprises a nonlinear function.

1 22. The computing unit of claim 21, wherein the nonlinear function
2 comprises a sigmoid function.

1 23. The computing unit of claim 21, wherein the nonlinear function
2 comprises a limiter function.

1 24. The computing unit of claim 12, further comprising a weight matrix,
2 wherein the weight matrix comprises the middle layer and output layer weights.

1 25. The computing unit of claim 12, wherein a node in the middle layer is
2 configured to process the input signal using middle layer weights by computing a dot product
3 of the middle layer weights and input signal for the connection to the node.

1 26. The computing unit of claim 25, wherein the node in the middle layer
2 is configured to process the dot product by applying the activation function to the dot
3 product.

1 27. The computing unit of claim 12, wherein a node in the output layer is
2 configured to process the middle layer signal using the output layer weights by computing a
3 dot product of the output layer weights and middle layer signal for the connections to the
4 node.

1 28. The computing unit of claim 27, wherein the node in the output layer is
2 configured to process the dot product by applying the activation function to the dot product.

1 29. The computing unit of claim 12, wherein the weights determine the
2 connection of nodes.

1 30. The computing unit of claim 12, wherein the computing machine
2 comprises an integrated circuit.

1 31. A method for performing a plurality of computing operations with a
2 computing unit using a weight matrix and an activation function, the computing unit
3 comprising a hardware structure that implements networked nodes, wherein nodes in the
4 networked nodes are related by a network of connections between the nodes, wherein the
5 weight matrix comprises weights corresponding to the connections and the activation
6 function specifies a function for the nodes in the networked nodes, the method comprising:
7 receiving an instruction that is applied to an input signal at the computing unit,
8 wherein the instruction includes the weight matrix and the activation function, the weight
9 matrix and activation function corresponding to a computing operation; and
10 mapping the input signal through the network of connections in the networked
11 nodes using the corresponding weights of the weight matrix for the connections and function
12 of the activation function for the nodes to generate an output signal, wherein the output signal
13 is a result of the computing operation determined by the weight matrix and activation
14 function.

1 32. The method of claim 31, wherein the networked nodes form a multi-
2 layer perceptron network.

1 33. The method of claim 32, wherein the multi-layer perceptron network is
2 a three layer perceptron network.

1 34. The method of claim 32, wherein mapping the input signal through the
2 network of connections in the networked nodes using the corresponding weights of the
3 weight matrix for the connections comprises computing a dot product for a node, wherein the
4 dot product is a computation of values of nodes connected to the node and the corresponding
5 weights for the connections to the node.

1 35. A method for performing a plurality of computing operations with a
2 computing unit using a weight matrix and an activation function, the computing unit
3 comprising a hardware structure that implements networked nodes, wherein nodes in the
4 networked nodes are related by a network of connections between the nodes, wherein the
5 weight matrix comprises weights corresponding to the connections and the activation
6 function specifies a function for the nodes in the networked nodes, the method comprising:
7 receiving an input signal at an input layer in the networked nodes;
8 sending the input signal to one or more nodes in a middle layer that are related
9 by connections with the input layer;
10 receiving middle layer weights for the connections between the input layer and
11 middle layer from the weight matrix;
12 processing the input signal using the middle layer weights and the function of
13 the activation function to generate a middle layer signal;
14 sending the middle layer signal to one or more nodes in an output layer that
15 are related by connections with the middle layer;
16 receiving output layer weights for the connections between the middle layer
17 and output layer from the weight matrix; and
18 generating an output signal by processing the middle layer signal using the
19 weights and the function of the activation function.

1 36. The method of claim 35, wherein processing the input signal using the
2 middle layer weights comprises computing a dot product for a node, wherein the dot product
3 is between values of nodes connected to the node and middle layer weights for the
4 connections to the node.

1 37. The method of claim 36, wherein processing the input signal using the
2 function of the activation function comprises computing the function of the dot product.

1 38. The method of claim 35, wherein processing the middle layer signal
2 using the middle layer weights comprises computing a dot product for a node, wherein the
3 dot product is between the middle layer signal and output layer weights for the connections to
4 the node.

1 39. The method of claim 38, wherein processing the middle layer signal
2 using the function of the activation function comprises computing the function of dot product.

1 40. A universal computing unit in a computing machine, wherein the
2 computing machine maps an input signal to an output signal using the universal computing
3 unit, the universal computing unit comprising:

4 a first layer configured to receive the input signal;

5 a second layer coupled to the first layer, the second layer comprising one or
6 more multiply-accumulate (MAC) units and one or more activation function modules,
7 wherein the one or more MAC units are configured to receive the input signal and second
8 layer weights from a weight matrix and calculate one or more dot products of the received
9 second layer weights and input signal, wherein the one or more activation function modules
10 are configured to calculate a function of the one or more dot products of the received second
11 layer weights and input signal to generate a second layer signal; and

12 a third layer coupled to the second layer, the third layer comprising one or
13 more MAC units and one or more activation function modules, wherein the one or more
14 MAC units are configured to receive the second layer signal and third layer weights from the
15 weight matrix and calculate one or more dot product of the received third layer weights and
16 second layer signal, wherein the one or more activation function modules are configured to
17 calculate a function of the one or more dot products of the received third layer weights and
18 second layer signal to generate the output signal.

1 41. The universal computing unit of claim 40, further comprising a weight
2 matrix module configured to send the second and third layer weights to the one or more MAC
3 units of the second and third layers.

1 42. The universal computing unit of claim 40, further comprising an
2 activation function module configured to send the function to the one or more activation
3 function modules.

1 43. The universal computing unit of claim 40, wherein the first layer
2 comprises a multiplexer.

1 44. The universal computing unit of claim 40, wherein the second layer
2 comprises a multiplexer configured to receive the second layer signal and send the second
3 layer signal to the one or more MACs of the third layer.

1 45. The universal computing unit of claim 40, wherein the third layer
2 comprises a multiplexer configured to provide the output signal.

1 46. A method for performing a plurality of computing operations with one
2 or more universal computing units, the one or more universal computing units being part of a
3 network that couples the one or more universal computing units to one or more computing
4 units, the method comprising:

5 receiving routing coefficients that specify connectivity information for the one
6 or more universal computing units and one or more computing units in the network, wherein
7 the routing coefficients replace a programming instruction stream by a data coefficient
8 stream;

9 connecting the one or more universal computing units and one or more
10 computing units in the network based on the routing coefficients;

11 receiving an instruction through the connected network comprising a weight
12 matrix and a selection of an activation function, wherein the weight matrix and selection of
13 the activation function comprise a set of operation-coefficients that define a desired
14 computing operation in the plurality of computing operations;

15 receiving an input data stream through the connected network; and
16 mapping an output data stream for the input data stream using the connected
17 one or more universal computing units and one or more computing units and the set of
18 operation-coefficients, the output data stream being a result of the defined desired computing
19 operation.

1 47. A system for performing a plurality of computing operations using one
2 or more universal computing units, the system comprising:
3 one or more computing units, wherein the one or more computing units form a
4 network with the one or more universal computing units, wherein the network is configured
5 to receive routing coefficients that specify connectivity information for the network, wherein
6 the routing coefficients replace a programming instruction stream by a data coefficient
7 stream,
8 wherein the one or more universal computing units receive an instruction
9 through the connected network, the instruction comprising a weight matrix and a selection of
10 an activation function, wherein the weight matrix and selection of the activation function
11 comprise a set of operation-coefficients that define a desired computing operation in the
12 plurality of computing operations, wherein the operation-coefficients replace a programming
13 instruction stream by a data coefficient stream.

11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65
66
67
68
69
70
71
72
73
74
75
76
77
78
79
80
81
82
83
84
85
86
87
88
89
90
91
92
93
94
95
96
97
98
99
100